

(19)日本国特許庁(JP)

(12)公表特許公報(A)

(11)特許出願公表番号

特表2023-537193  
(P2023-537193A)

(43)公表日 令和5年8月31日(2023.8.31)

(51)Int.Cl. F I  
G 0 6 F 16/906 (2019.01) G 0 6 F 16/906

テーマコード(参考)  
5 B 1 7 5

審査請求 有 予備審査請求 未請求 (全 15 頁)

(21)出願番号 特願2022-578769(P2022-578769)  
(86)(22)出願日 令和4年3月15日(2022.3.15)  
(85)翻訳文提出日 令和5年2月13日(2023.2.13)  
(86)国際出願番号 PCT/IB2022/052333  
(87)国際公開番号 WO2022/269368  
(87)国際公開日 令和4年12月29日(2022.12.29)  
(31)優先権主張番号 202141028706  
(32)優先日 令和3年6月25日(2021.6.25)  
(33)優先権主張国・地域又は機関  
インド(IN)

(71)出願人 520412811  
エルアンドティー テクノロジー サービ  
シズ リミテッド  
インド国, チェンナイ 6 0 0 0 8 9, ラ  
ーマプラーム, マウント プーナマリー  
ロード, ディーエルエフ アイティー エ  
スイーゼット パーク, ブロック 3, 2  
番 フロア, 1 / 1 2 4  
(74)代理人 110002365  
弁理士法人サンネクスト国際特許事務所  
(72)発明者  
ダス, イシタ  
インド 7 0 0 0 0 8 ウエストベンガル  
, コルカタ, イーストパーク, バ  
リシャ, カリパダ ムカルジー ロード  
, 3 3 9 / ジー / 1 7

最終頁に続く

(54)【発明の名称】 クラスタを表現するためにサンプルを選択する方法およびシステム

(57)【要約】

クラスタを表現するためにサンプルを選択する方法が開示される。方法は、最適化デバイスによって、1つまたは複数のクラスタを受信することを含んでもよい。1つまたは複数のクラスタのそれぞれは、複数のサンプルを含む。方法は、1つまたは複数のクラスタのそれぞれから選択されるサンプルの数のカウントを決定してもよく、1つまたは複数のクラスタのそれぞれについて、アレイベースの距離行列を生成してもよい。方法は、クラスタの複数のサンプルを、クラスタ内の複数のサンプルのばらつき具合に基づいて、ソートしてもよい。ソートすることは、1つまたは複数のクラスタのそれぞれについてのアレイベースの距離行列を使用して実施してもよい。さらに、方法は、クラスタを表現するために、サンプルの数の決定されたカウントを、複数のクラスタのそれぞれの、ソートされた複数のサンプルから選択してもよい。

【選択図】 図 3

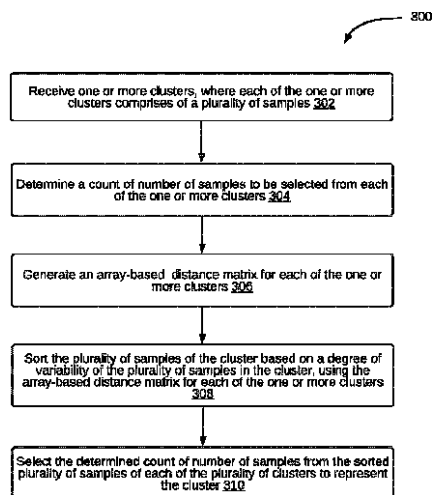


FIG. 3

## 【特許請求の範囲】

## 【請求項1】

クラスタを表現するためにサンプルを選択する方法であって、最適化デバイスによって、それぞれが複数のサンプルを備える1つまたは複数のクラスタを受信し、前記最適化デバイスによって、前記1つまたは複数のクラスタのそれぞれから選択されるサンプルの数のカウントを決定し、前記最適化デバイスによって、前記1つまたは複数のクラスタのそれぞれについて、アレイベースの距離行列を生成し、前記最適化デバイスによって、前記クラスタの前記複数のサンプルを、前記クラスタ内の前記複数のサンプルのばらつきの度合いに基づいて、前記1つまたは複数のクラスタのそれぞれについての前記アレイベースの距離行列を使用して、ソートし、前記最適化デバイスによって、前記クラスタを表現するために、サンプルの数の前記決定されたカウントを、前記複数のクラスタのそれぞれの、前記ソートされた複数のサンプルから選択する、方法。

10

## 【請求項2】

前記サンプルの数の前記選択された決定されたカウントがしきい値未満である場合、前記サンプルの数の所定のカウントを、前記1つまたは複数のクラスタのそれぞれから選択し、前記しきい値が、各データセットに特有であり、クラスタのサイズを前記データセットと比較することによって決定される、請求項1に記載の方法。

20

## 【請求項3】

前記1つまたは複数のクラスタのそれぞれから選択されるサンプルの数の前記カウントが、前記1つまたは複数のクラスタのそれぞれについてのサイズ、ばらつき、およびクラスタ確率のうちの少なくとも1つに基づいて、層化サンプリング技法を使用して、決定される、請求項1に記載の方法。

## 【請求項4】

前記クラスタ確率が、機械学習（ML）モデルを使用して決定され、前記MLモデルが、前記クラスタの前記複数のサンプルを分類する、請求項3に記載の方法。

30

## 【請求項5】

下記を行なうように構成された1つまたは複数のコンピューティングデバイスを備える、最適化デバイスによって、それぞれが複数のサンプルを備える1つまたは複数のクラスタを受信する、前記最適化デバイスによって、前記1つまたは複数のクラスタのそれぞれから選択されるサンプルの数のカウントを決定する、前記最適化デバイスによって、前記1つまたは複数のクラスタのそれぞれについて、アレイベースの距離行列を生成する、前記最適化デバイスによって、前記クラスタの前記複数のサンプルを、前記クラスタ内の前記複数のサンプルのばらつきの度合いに基づいて、前記1つまたは複数のクラスタのそれぞれについての前記アレイベースの距離行列を使用して、ソートする、および、前記最適化デバイスによって、前記クラスタを表現するために、サンプルの数の前記決定されたカウントを、前記複数のクラスタのそれぞれの、前記ソートされた複数のサンプルから選択する、システム。

40

## 【請求項6】

前記サンプルの数の前記選択された決定されたカウントがしきい値未満である場合、前記サンプルの数の所定のカウントを、前記1つまたは複数のクラスタのそれぞれから選択し

50

、前記しきい値が、各データセットに特有であり、クラスタのサイズを前記データセットと比較することによって決定される、請求項5に記載のシステム。

【請求項7】

前記1つまたは複数のクラスタのそれぞれから選択されるサンプルの数の前記カウントが、1つまたは複数のクラスタのそれぞれについてのサイズ、ばらつき、およびクラスタ確率のうち少なくとも1つに基づいて、層化サンプリング技法を使用して、決定される、請求項5に記載のシステム。

【請求項8】

前記クラスタ確率が、機械学習（ML）モデルを使用して決定され、前記MLモデルが、前記クラスタの前記複数のサンプルを分類する、請求項7に記載のシステム。

【発明の詳細な説明】

【技術分野】

【0001】

本開示は、一般に、データセットのサイズを小さくすることに関し、より詳細には、データセットのサイズを小さくするためのクラスタを表現するために複数のサンプルを選択することに関する。

【背景技術】

【0002】

情報爆発を伴うビッグデータの時代では、データ処理に関連するビジネス要件が、日々高まっており、処理対象のデータサンプルは、益々複雑化している。データのクラスタリングは、システムのエンジニアリングおよび計算科学において極めて重要な分野であり、教師無し学習を使用して、ラベル付けされたデータを伴わずに機械学習モデルを訓練する。クラスタリングアルゴリズムは、自身の一意な特徴を有するいくつかのクラスタに、データを分割する。時には、これらのクラスタ自身が、膨大な数のサンプルを有する。ラベル付けされていないデータが利用可能となるのは、高い次元、また線形に分離不可能なデータ空間であり、そのことにより、機械学習モデルの処理および訓練の間、大きなメモリチャックおよび時間を消費する。

【発明の概要】

【0003】

したがって、当分野では、機械学習モデルの効果的でリソース効率の良い訓練のために、クラスタ内のサンプルの総数のカウントを低減することによって、データセットのサイズを小さくするための方法およびシステムを提供する必要性がある。

【0004】

クラスタを表現するためにサンプルを選択する方法が開示される。方法は、最適化デバイスによって、1つまたは複数のクラスタを受信することを含んでもよい。1つまたは複数のクラスタのそれぞれは、複数のサンプルを含む。方法は、1つまたは複数のクラスタのそれぞれから選択されるサンプルの数のカウントを決定してもよく、1つまたは複数のクラスタのそれぞれについて、アレイベースの距離行列を生成してもよい。方法は、クラスタの複数のサンプルを、クラスタ内の複数のサンプルのばらつきの度合いに基づいて、ソートしてもよい。ソートすることは、1つまたは複数のクラスタのそれぞれについてのアレイベースの距離行列を使用して実施してもよい。さらに、方法は、クラスタを表現するために、サンプルの数の決定されたカウントを、複数のクラスタのそれぞれの、ソートされた複数のサンプルから選択してもよい。

【0005】

本開示に組み込まれ、その一部を成す添付の図面は、例示の実施形態を図示しており、説明と併せて、開示される原理を説明するよう機能する。

【図面の簡単な説明】

【0006】

10

20

30

40

50

【図1】本開示の実施形態による、1つまたは複数のクラスタから複数のデータサンプルを選択するためのプロセス図である。

【図2】本開示のいくつかの実施形態による、1つまたは複数のクラスタから複数のデータサンプルをソートして選択するためのプロセス図である。

【図3】本開示のいくつかの実施形態による、クラスタを表現するためにサンプルを選択する方法のフローチャートである。

【発明を実施するための形態】

【0007】

例示的な実施形態を、添付の図面を参照して説明する。同一または類似の部分参照するために、好都合であればいつでも、図面を通じて同一の参照符号を使用する。本明細書では開示される原理の例および特徴が説明されるが、開示される実施形態の思想および範囲から逸脱することなく、修正、適合、および他の実装が可能である。以下の詳細な説明は、単なる例として考えられ、真の範囲および思想は、後に続く特許請求の範囲によって示されることが意図されている。追加的な例示の実施形態を以降に列挙する。

10

【0008】

理解されるように、クラスタリングアルゴリズムは、自身の一意的特徴を有するいくつかのクラスタに、データを分割する。時には、これらのクラスタ自身が、膨大な数のサンプルを有する。本開示は、クラスタに固有なばらつき、性質をカバーする、限られた数のサンプルを用いてクラスタを表現することができる解決策を提供する。このやり方で、アルゴリズムは、さらなるプロセスのためにデータセット全体を使用するために依存性を低減し、それにより、大きなデータセットを用いて作業する際のメモリおよび時間の複雑さを制限する。アルゴリズムはまた、クラスタのサイズが小さい場合に、ユーザが必要なサンプル数をクラスタから選択できるようにする、柔軟なものである。

20

【0009】

さらには、プロセスは、均質なクラスタからでも、一意的サンプルを選択することができるよう確実にする。

【0010】

図1を参照すると、本開示の実施形態による、1つまたは複数のクラスタから複数のデータサンプルを選択するためのプロセス100が図示される。

30

【0011】

ステップ102では、データセットは、1つまたは複数の異なるクラスタにクラスタ化されてもよい。クラスタリングは、同じように見え、類似の特徴を有するデータセットが、特定のクラスタ内に一緒に保持されることを確実にするよう実施されてもよい。

【0012】

ステップ104では、複数のデータサンプルのうち、いくつかのデータサンプルを、1つまたは複数の異なる作成済クラスタのうち1つのクラスタから選択することができるかを決定することができる。

【0013】

いくつかのデータサンプルが選択されるかを決定するために、最適な割り当てのための、層化サンプリングメカニズムが使用されてもよい。層化サンプリングメカニズムは、複数のデータサンプルを考慮することができる。複数のデータサンプルのそれぞれは、均質なグループ（つまり、類似の特徴を有する複数のデータサンプルのそれぞれを一緒に格納することができる、クラスタ）に分割することができる。決定は、複数の類似して見えるサンプルの中から、いくつかのサンプルが選択され得るかに関する場合がある。ステップ106では、データサンプルのうち、どれを、1つまたは複数の異なるクラスタから選択することができるかを決定することができる。層化サンプリングメカニズムは、特定の均質なデータグループのうち1つを選択してもよく、特定の計算に基づいて1つまたは複数のデータサンプルをランダムに選択してもよい。

40

【0014】

ステップ108では、所与のクラスタについて、下で言及する式を用いて、 $N_i$ 個のサン

50

プルをクラスタから選択することができる：

【数1】

$$n_i = \frac{w_i S_i}{\sqrt{c_i} \sum_{i=1}^k \sqrt{c_i} w_i S_i} \cdot \frac{C_0}{C_0} \quad \dots \quad \text{式(1)}$$

ただし、 $w_i$ は、 $i$ 番目のクラスタ内に存在するデータサンプルの数であり、  
 $S_i$ は、クラスタ内のデータサンプルのばらつきであり、  
 $C_i$ は、平均のクラスタ確率であり、  
 $C_0$ は、定数である。

10

【0015】

式(1)は、式(2)のように、簡単にすることができる。

【数2】

$$n = \sum_{i=1}^k n_i \quad \dots \quad \text{式(2)}$$

20

【0016】

式(2)では、クラスタのサイズとばらつきを考慮してもよい。

【0017】

クラスタ内に存在する複数のデータサンプル同士を区別するために機械学習モデルが利用可能であり、また所与のセットの特徴に基づいて複数のデータサンプル同士の区別の正確なレベルを予測することができる事例では、機械学習モデルから導かれる学習を利用して、クラスタ確率を考慮することによって、いくつのデータサンプルをクラスタから選択することができるかを決定することができる。区別のレベルに関する学習が、利用可能な機械学習モデルのいずれからも利用可能ではない場合、クラスタ確率パラメータは、データサンプルの数の選択を決定するための重み付けが、あまり与えられない場合がある。

30

【0018】

ステップ106では、どのデータサンプルが選択されるかに関する決定が実施されてもよい。一般的に、データサンプルは、利用可能なランダム選択メカニズムのいずれかに基づいて、ランダムに選択され得る。しかしながら、選択されたデータサンプルがクラスタ全体を表現するべく、クラスタのばらつきの度合いを最大化するために、距離ベースの選択メカニズムが、ステップ110で利用されてもよい。距離ベースの選択メカニズムでは、クラスタ内に存在するアレイベースで最適化された距離が利用されてもよい。距離行列は、例えば、ユークリッドベースの距離行列またはマンハッタンベースの距離行列であつてもよい。さらに、データサンプルは、その距離に基づいて、つまり、ばらつきの最大化に基づいてソートされてもよい。

40

【0019】

ステップ116では、1つまたは複数のクラスタの、クラスタのそれぞれにおいて、式(2)を使用して、「 $n_i$ 」個のデータサンプルが選択されてもよい。さらに、「 $n_i$ 」個のデータサンプルを選択する手順は、データセットの1つまたは複数のクラスタのクラスタすべてについて繰り返すことができる。具体的なシナリオでは、クラスタから選択された、ある数のデータサンプルが最小限である場合、プロセス100は、1つまたは複数のクラスタのそれぞれから、サンプルの数の所定のカウンタを選択してもよい。これは、サンプルの数の選択された決定されたカウンタが、しきい値未満である場合に行なわれ得る

50

。ステップ118では、合計「n」個のデータサンプルが、1つまたは複数のクラスタのそれぞれから選択され、それにより、データセットのサイズを小さくする。

【0020】

次に図2を参照すると、本開示のいくつかの実施形態による、1つまたは複数のクラスタからの複数のデータサンプルをソートして選択するためのプロセス200が図示される。

【0021】

ステップ204では、1つまたは複数のクラスタのうち、1つのクラスタから第1のデータサンプルが選択されてもよい。ステップ206では、第1の選択されたサンプルから最も遠い、第2のデータサンプルが選択されてもよい。ステップ208では、第1のデータサンプルと第2のサンプルは、データセット内に維持されてもよい。ステップ210では、第3のデータサンプルが選択されてもよい。第3のデータサンプルの選択は、ステップ216において、あるメカニズムにしたがって実施され、データセットの外部から、例えば第3のデータサンプルから、ランダムなサンプルが選択される。データセットの外部のデータサンプルに関して、データセットのデータサンプルの距離を、決定することができる。例えば、第3のデータサンプルの距離は、データセットの第1のデータサンプルに関して「d13」として、およびデータセットの第2のデータサンプルに関して「d23」として決定されてもよい。

10

【0022】

距離「d13」と距離「d23」の距離のうち、小さいほうを選択することができる。小さいほうの距離とは、例えば図2で図示されるような「d13」である。さらに、データセットの外部に存在するデータサンプルのすべてについて、距離をチェックして最小距離を選択する、上で言及したステップが決定されてもよい。例えば、データセットの外部の別のデータサンプルは、第4のデータサンプルであってもよく、決定される距離は、例えば、データセットの第1のデータサンプルから第4のデータサンプルまで、「d14」として、およびデータセットの第2のデータサンプルからの「d24」としてであってもよい。

20

【0023】

続いて、決定される距離の最小を、例えば「d13」および「d24」として決定することができる。次に、最小の決定された距離からの最大距離、例えば「d13」が選択される。最終的に、データサンプル、例えば最大距離に対応する第3のサンプル、例えば「d13」が選択されてもよく、データセットに挿入されてもよい。

30

【0024】

層化サンプリングメカニズムを使用することによって、および、ばらつきの最大化に基づいてクラスタ内でデータサンプルをソートすることによって、ステップ212では、「n」個のサンプルを、1つまたは複数のクラスタのクラスタから選択することができ、それによって、選択されたサンプルは一意であり、そのクラスタのばらつき全体をカバーすることができる。

【0025】

ステップ214では、上述のステップ204～212は、1つまたは複数のクラスタの、クラスタのそれぞれについて繰り返され、データセットの性質を保った新しい縮小データセットを作成することができる。

40

【0026】

例示の実施形態では、手書きのアルファベットのデータセットがあると仮定する。データセットでは、アルファベット「a」は、様々なユーザによって、イタリック体、ボールド体、異なるフォントサイズで、または筆記体など、異なった形で書かれ得る。さらに、アルファベットは、アルファベットが「a」、「b」、「c」などのアルファベットのカテゴリに存在するかどうかに基づいてクラスタ化され得る。アルファベット「a」のクラスタからの複数のデータサンプルを検討する。アルファベット「a」のクラスタでは、イタリック体の「a」は、数が少ないものとする。例えば、50Kデータサンプルを有するデータセット中、5Kデータサンプルだけが、アルファベット「a」のイタリック体を表現

50

する場合があります、故に50Kデータサンプル内で、イタリック体の「a」の一意性を表現し得る。一意なイタリック体「a」は、データサンプルを配列してソートするために使用されてもよい。したがって、50Kサンプルから、5Kサンプルが、アルファベット「a」のイタリック体を表現するために使用することができると結論付けることができる。さらに、これらの表現される5Kサンプルのうち、どれが選ばれるかは、クラスタ内でのデータサンプルのばらつきの最大化に基づいたソーティングメカニズムによって決定されてもよい。

#### 【0027】

次に図3を参照すると、実施形態による、クラスタを表現するためにサンプルを選択する方法300のフローチャートが図示される。ステップ302では、1つまたは複数のクラスタが受信される。1つまたは複数のクラスタのそれぞれは、複数のサンプルを含んでもよい。

10

#### 【0028】

ステップ304では、1つまたは複数のクラスタのそれぞれから選択されるサンプルの数のカウントが決定されてもよい。1つまたは複数のクラスタのそれぞれから選択されるサンプルの数のカウントは、1つまたは複数のクラスタのそれぞれについてのサイズ、ばらつき、およびクラスタ確率のうち少なくとも1つに基づいて、層化サンプリング技法を使用して、決定されてもよい。クラスタ確率は、機械学習(ML)モデルを使用して決定されてもよく、MLモデルは、クラスタの複数のサンプルを分類する。訓練無しMLモデルの事例では、各クラスタは、等しい確率を割振られることに留意されたい。

20

#### 【0029】

ステップ306では、1つまたは複数のクラスタのそれぞれについて、アレイベースの距離行列が生成されてもよい。例えば、アレイベースの距離行列は、ユークリッド距離行列であってもよい。ステップ308では、クラスタの複数のサンプルは、クラスタ内の複数のサンプルのばらつきの度合いに基づいて、1つまたは複数のクラスタのそれぞれについてのアレイベースの距離行列を使用して、ソートすることができる。

#### 【0030】

ステップ310では、クラスタを表現するために、サンプルの数の決定されたカウントを、複数のクラスタのそれぞれの、ソートされた複数のサンプルから選択してもよい。加えて、いくつかの実施形態では、サンプルの数の選択された決定されたカウントがしきい値未満である場合、サンプルの数の所定のカウントは、1つまたは複数のクラスタのそれぞれから選択されてもよい。

30

#### 【0031】

1つまたは複数のコンピュータ可読記憶媒体は、本開示と一貫性のある実施形態を実装する際に利用される。コンピュータ可読記憶媒体とは、プロセッサによって可読な情報またはデータを記憶することができる、あらゆるタイプの物理メモリを称する。故に、コンピュータ可読記憶媒体は、本明細書で説明される実施形態と一貫性のあるステップまたは工程をプロセッサに実行させるための命令を含む、1つまたは複数のプロセッサによる実行のための命令を記憶する。「コンピュータ可読媒体」という用語は、有形物を含み、搬送波および一過性の信号を含まない、つまり非一時的であると理解されるべきである。例としては、ランダムアクセスメモリ(RAM)、読み取り専用メモリ(ROM)、揮発性メモリ、非揮発性メモリ、ハードドライブ、CD-ROM、DVD、フラッシュドライブ、ディスク、およびあらゆる他の既知の物理的な記憶媒体が挙げられる。

40

#### 【0032】

明確にするために、上述のことは、本開示の実施形態を、異なる機能的なユニットおよびプロセッサに関して説明したものであることを諒解されたい。しかしながら、異なる機能的なユニット、プロセッサ、またはドメイン間での機能性のあらゆる好適な分散が、本開示から逸脱することなく使用され得ることが明らかとなろう。

#### 【0033】

例えば、別個のプロセッサまたはコントローラによって実施されるよう図示される機能性

50

は、同一のプロセッサまたはコントローラによって実施されてもよい。したがって、特定の機能的なユニットへの参照は、厳密な論理的または物理的な構造または編成を示すのではなく、説明される機能性を提供するための好適な手段への参照として考えられるに過ぎない。

【0034】

いくつかの実施形態に関連して本開示を説明したが、本開示は、本明細書で説明した特定の形態に限定されるよう意図されていない。そうではなく、本開示の範囲は、特許請求の範囲によってのみ制限される。加えて、特徴は、特定の実施形態に関連して説明されるように見えるかもしれないが、当業者であれば、説明される実施形態の様々な特徴は、本開示にしたがって組み合わせることができることを認識されよう。

10

【0035】

さらには、個々に列挙されているが、複数の手段、要素、またはプロセスのステップは、例えば単一のユニットまたはプロセッサによって実装されてもよい。加えて、個々の特徴が、異なる請求項に含まれる場合があるが、これらは、可能であれば有利に組み合わせられてもよく、異なる請求項に含まれることは、特徴の組み合わせが、実行可能ではない、および/または有利ではないということを意味するものではない。また、特徴が、請求項の1つのカテゴリに含まれることは、このカテゴリへの限定することを意味するものではなく、むしろ、適当であれば特徴は他の請求項カテゴリに等しく適用可能であり得る。

20

30

40

50

【図 1】

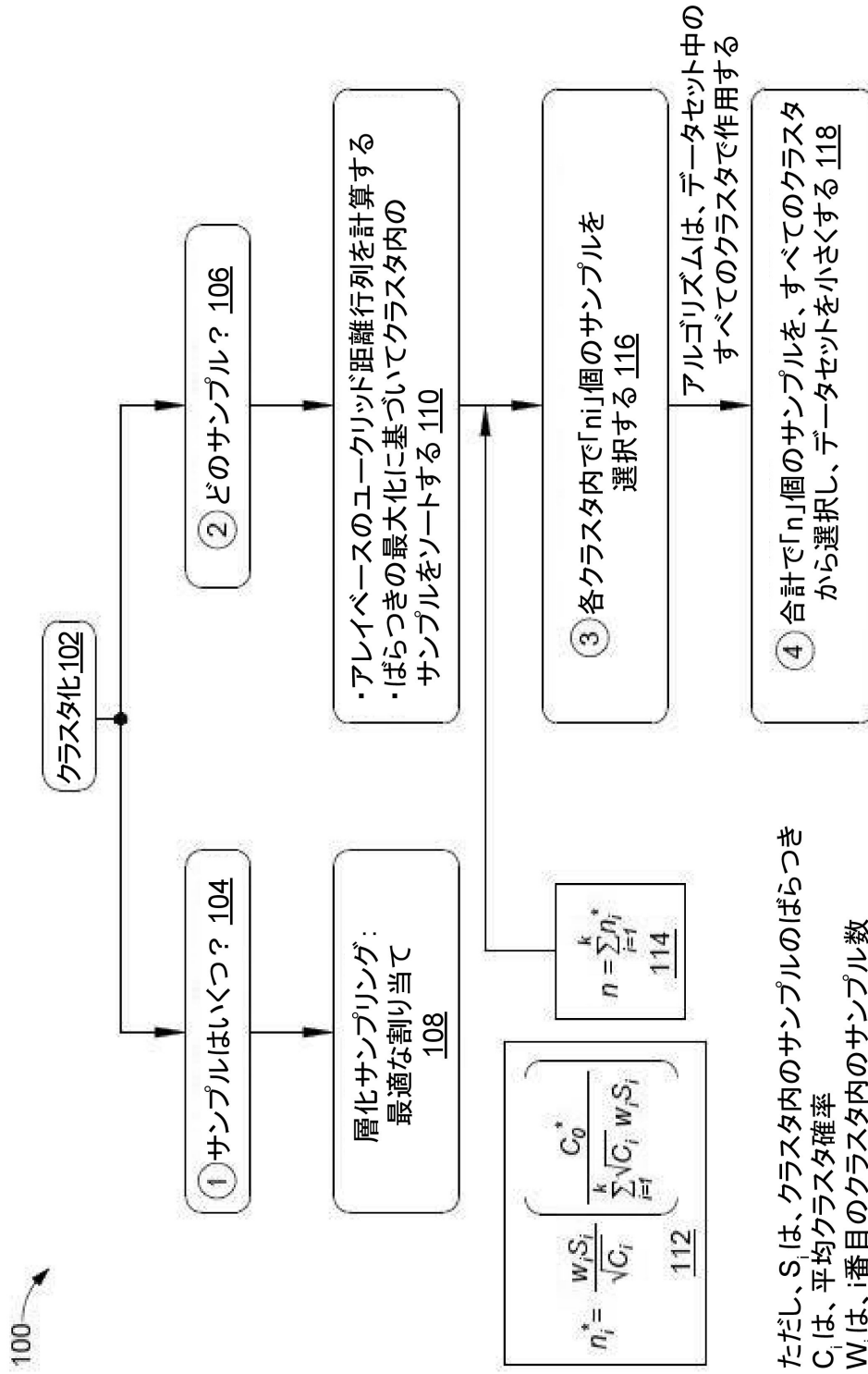


図 1

【図2】

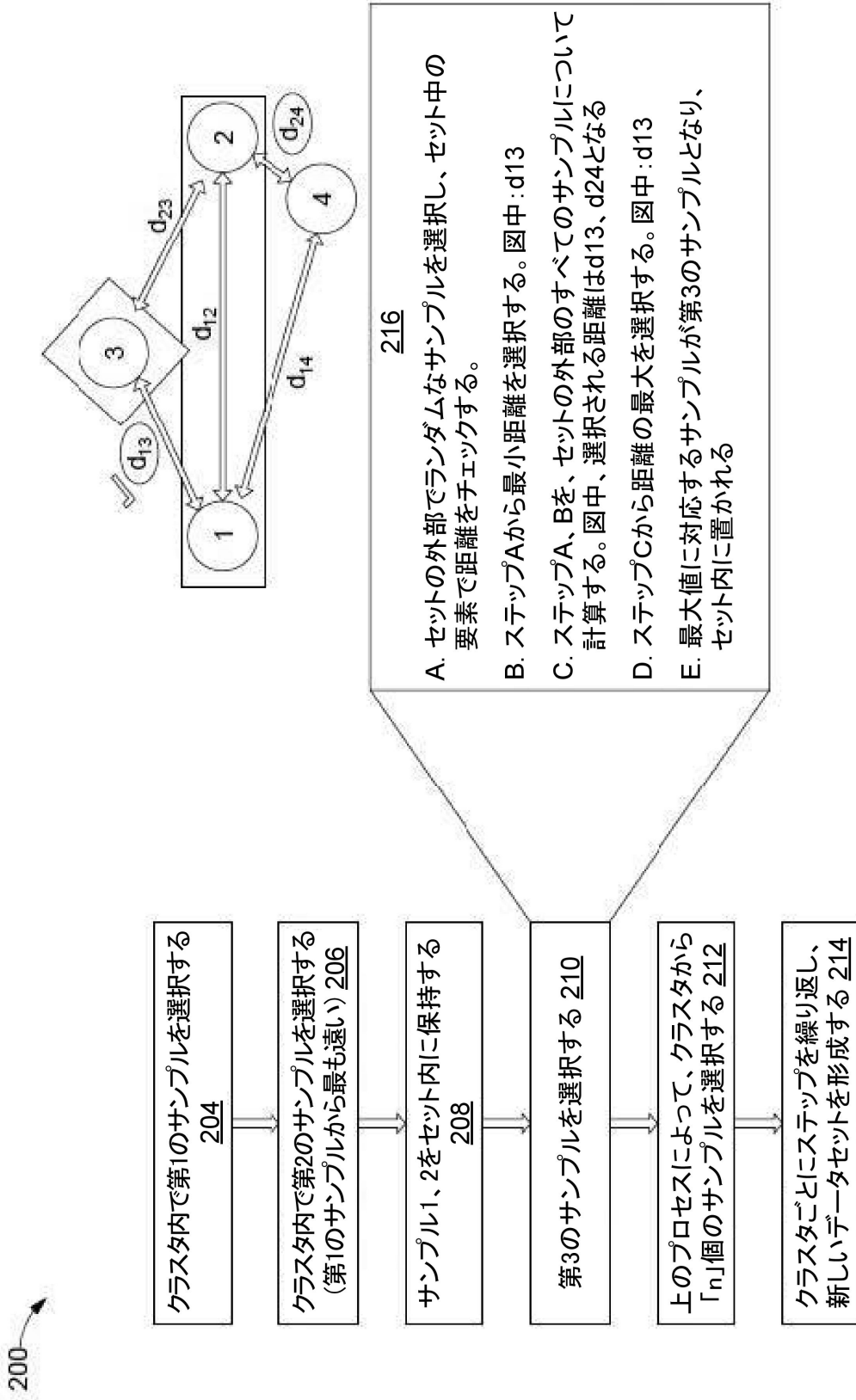


図2

10

20

30

40

50

【図 3】

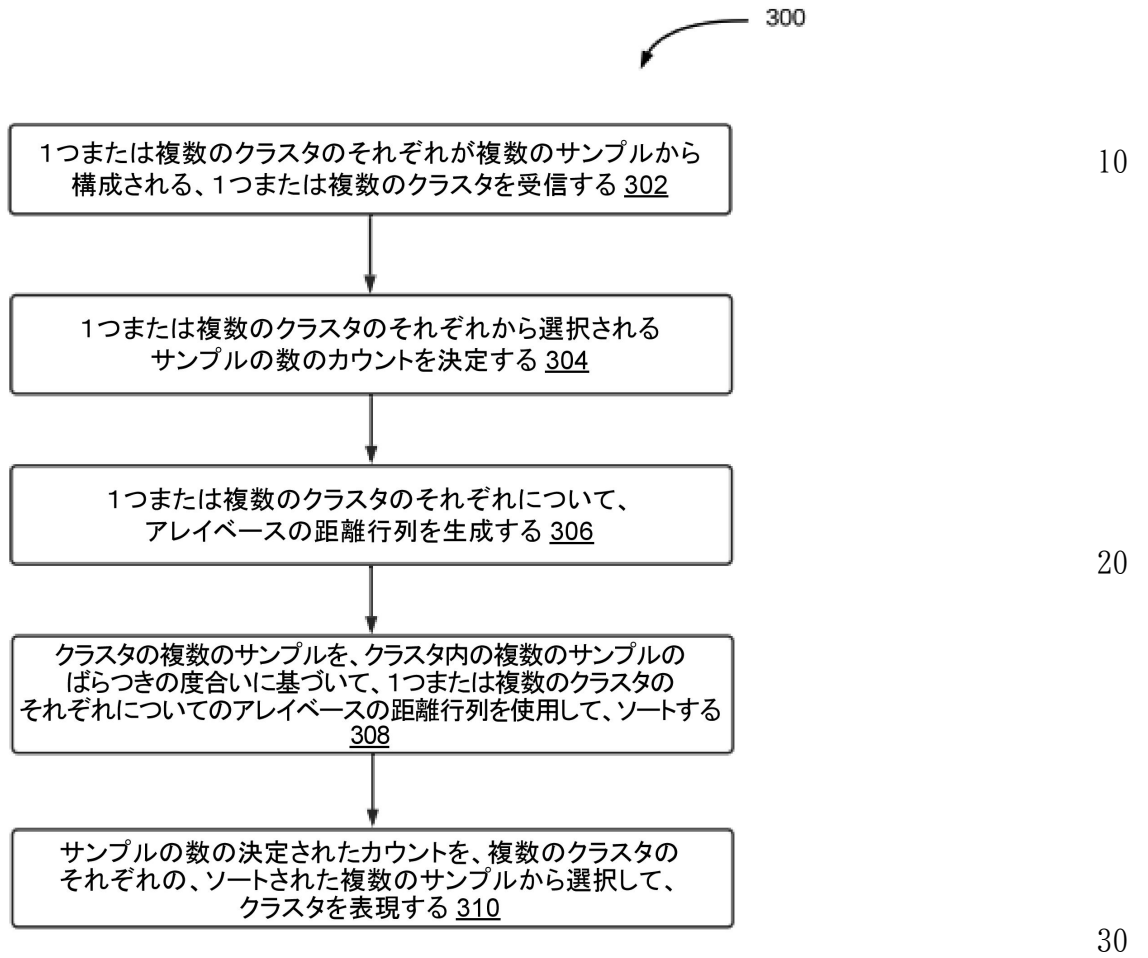


図 3

10

20

30

40

50

【手続補正書】

【提出日】令和4年8月16日(2022.8.16)

【手続補正1】

【補正対象書類名】特許請求の範囲

【補正対象項目名】全文

【補正方法】変更

【補正の内容】

【特許請求の範囲】

【請求項1】

クラスタを表現するためにサンプルを選択する方法であって、  
最適化デバイスによって、それぞれが複数のサンプルを備える1つまたは複数のクラスタを受信し(306)、  
前記最適化デバイスによって、前記1つまたは複数のクラスタのそれぞれについてのサイズ、ばらつき、およびクラスタ確率のうち少なくとも1つに基づいて、層化サンプリング技法を使用して、前記1つまたは複数のクラスタのそれぞれから選択されるサンプルの数のカウントを決定し(304)、  
前記最適化デバイスによって、前記1つまたは複数のクラスタのそれぞれについて、アレイベースの距離行列を生成し(306)、  
前記最適化デバイスによって、前記クラスタの前記複数のサンプルを、前記クラスタ内の前記複数のサンプルのばらつきの度合いに基づいて、前記1つまたは複数のクラスタのそれぞれについての前記アレイベースの距離行列を使用して、ソートし(308)、  
前記最適化デバイスによって、前記クラスタを表現するために、サンプルの数の前記決定されたカウントを(310)、前記複数のクラスタのそれぞれの、前記ソートされた複数のサンプルから選択する、  
方法。

10

20

【請求項2】

前記サンプルの数の前記選択された決定されたカウントがしきい値未満である場合、前記サンプルの数の所定のカウントを、前記1つまたは複数のクラスタのそれぞれから選択し、前記しきい値が、各データセットに特有であり、クラスタのサイズを前記データセットと比較することによって決定される、  
請求項1に記載の方法。

30

【請求項3】

前記クラスタ確率が、機械学習(ML)モデルを使用して決定され、前記MLモデルが、前記クラスタの前記複数のサンプルを分類する、  
請求項1に記載の方法。

【請求項4】

下記を行なうように構成された1つまたは複数のコンピューティングデバイスを備える、最適化デバイスによって、それぞれが複数のサンプルを備える1つまたは複数のクラスタを受信する、  
前記最適化デバイスによって、前記1つまたは複数のクラスタのそれぞれについてのサイズ、ばらつき、およびクラスタ確率のうち少なくとも1つに基づいて、層化サンプリング技法を使用して、前記1つまたは複数のクラスタのそれぞれから選択されるサンプルの数のカウントを決定する、  
前記最適化デバイスによって、前記1つまたは複数のクラスタのそれぞれについて、アレイベースの距離行列を生成する、  
前記最適化デバイスによって、前記クラスタの前記複数のサンプルを、前記クラスタ内の前記複数のサンプルのばらつきの度合いに基づいて、前記1つまたは複数のクラスタのそれぞれについての前記アレイベースの距離行列を使用して、ソートする、および、  
前記最適化デバイスによって、前記クラスタを表現するために、サンプルの数の前記決定されたカウントを、前記複数のクラスタのそれぞれの、前記ソートされた複数のサンプル

40

50

から選択する、  
システム。

【請求項5】

前記サンプルの数の前記選択された決定されたカウントがしきい値未満である場合、前記サンプルの数の所定のカウントを、前記1つまたは複数のクラスタのそれぞれから選択し、前記しきい値が、各データセットに特有であり、クラスタのサイズを前記データセットと比較することによって決定される、  
請求項4に記載のシステム。

【請求項6】

前記クラスタ確率が、機械学習（ML）モデルを使用して決定され、前記MLモデルが、  
前記クラスタの前記複数のサンプルを分類する、  
請求項4に記載のシステム。

10

20

30

40

50

## 【国際調査報告】

INTERNATIONAL SEARCH REPORT		International application No. PCT/IB2022/052333
A. CLASSIFICATION OF SUBJECT MATTER G06N20/00 Version=2022.01		
According to International Patent Classification (IPC) or to both national classification and IPC		
B. FIELDS SEARCHED		
Minimum documentation searched (classification system followed by classification symbols) G06N		
Documentation searched other than minimum documentation to the extent that such documents are included in the fields searched		
Electronic data base consulted during the international search (name of data base and, where practicable, search terms used) Databases: IPO Internal Database, Patseer Keywords: Cluster, sample screening, re-grouping, representation		
C. DOCUMENTS CONSIDERED TO BE RELEVANT		
Category*	Citation of document, with indication, where appropriate, of the relevant passages	Relevant to claim No.
X	US20160180556A1 (DENG CHANG) (23-06-2016) 23 June 2016 Abstract, paragraphs [0003], [0007], [0009]-[0011], [0062]-[0063], [0065]-[0066], figure 4	1-8
X	----- CN107194430A (BEIJING SANKUAI ON LINE SCIENCE & TECHNOLOGY CO LTD) (22-09-2017) 22 September 2017 Whole Document	1-8
<input type="checkbox"/> Further documents are listed in the continuation of Box C. <input type="checkbox"/> See patent family annex.		
* Special categories of cited documents:      "T" later document published after the international filing date or priority date and not in conflict with the application but cited to understand the principle or theory underlying the invention "A" document defining the general state of the art which is not considered to be of particular relevance "D" document cited by the applicant in the international application      "X" document of particular relevance; the claimed invention cannot be considered novel or cannot be considered to involve an inventive step when the document is taken alone "E" earlier application or patent but published on or after the international filing date "L" document which may throw doubts on priority claim(s) or which is cited to establish the publication date of another citation or other special reason (as specified)      "Y" document of particular relevance; the claimed invention cannot be considered to involve an inventive step when the document is combined with one or more other such documents, such combination being obvious to a person skilled in the art "O" document referring to an oral disclosure, use, exhibition or other means "P" document published prior to the international filing date but later than the priority date claimed      "&" document member of the same patent family		
Date of the actual completion of the international search 21-06-2022		Date of mailing of the international search report 21-06-2022
Name and mailing address of the ISA/ Indian Patent Office Plot No.32, Sector 14, Dwarka, New Delhi-110075 Facsimile No.		Authorized officer Shri Ram Kaunaujia Telephone No. +91-1125300200

Form PCT/ISA/210 (second sheet) (July 2019)

10

20

30

40

50

---

 フロントページの続き

(81)指定国・地域 AP(BW, GH, GM, KE, LR, LS, MW, MZ, NA, RW, SD, SL, ST, SZ, TZ, UG, ZM, ZW), EA(AM, AZ, BY, KG, KZ, RU, TJ, TM), EP(AL, AT, BE, BG, CH, CY, CZ, DE, DK, EE, ES, FI, FR, GB, GR, HR, HU, IE, IS, IT, LT, LU, LV, MC, MK, MT, NL, NO, PL, PT, RO, RS, SE, SI, SK, SM, TR), OA(BF, BJ, CF, CG, CI, CM, GA, GN, GQ, GW, KM, ML, MR, NE, SN, TD, TG), AE, AG, AL, AM, AO, AT, AU, AZ, BA, BB, BG, BH, BN, BR, BW, BY, BZ, CA, CH, CL, CN, CO, CR, CU, CZ, DE, DJ, DK, DM, DO, DZ, EC, EE, EG, ES, FI, GB, GD, GE, GH, GM, GT, HN, HR, HU, ID, IL, IN, IR, IS, IT, JM, JO, JP, KE, KG, KH, KN, KP, KR, KW, KZ, LA, LC, LK, LR, LS, LU, LY, MA, MD, ME, MG, MK, MN, MW, MX, MY, MZ, NA, NG, NI, NO, NZ, OM, PA, PE, PG, PH, PL, PT, QA, RO, RS, RU, RW, SA, SC, SD, SE, SG, SK, SL, ST, SV, SY, TH, TJ, TM, TN, TR, TT, TZ, UA, UG, US, UZ, VC, VN, WS, ZA, ZM, ZW

(72)発明者 シン, マドウスダン  
 インド 560100 カルナータカ, バンガロール, エレクトロニック シティー-1,  
 ニーラドリー ロード, ファースト クロス, アジメラ ストーン パーク, ビー-603

(72)発明者 バララマン, ムリダル  
 インド 560037 カルナータカ, バンガロール, ドダネクンディ, チナパナハリ メ  
 インロード, エスブイエス パームズ 2, ビー 206

(72)発明者 デブナス, スカント  
 インド 382481 グジャラート, アフマダーバード, チャンドロディヤ, エヌアール  
 プラカシュ ナガー ソサエティ, パート-3, ガジュラジ ソサエティ, 177

(72)発明者 グプタ, ムリナル  
 インド 180002 ジャンムー カシミール, ジャンムー, タラブ ティロ, バサント  
 ビハール, 78-エー

Fターム(参考) 5B175 FA03 FB04